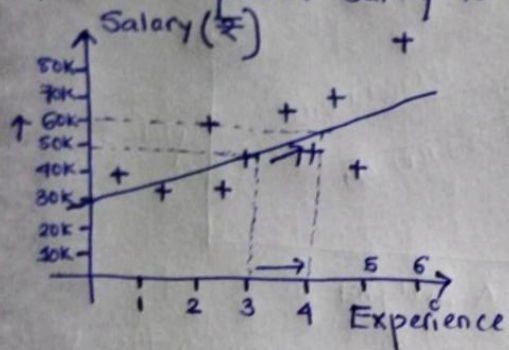


SIMPLE AND MULTIPLE LINEAR REGRESSION (OLS, Dummy Variable trap)

SIMPLE LINEAR REGRESSION, $y = b_0 + b_1 x_1$ $b_0 \rightarrow$ Constant
 $y \rightarrow$ Dependent Variable (DV), $x_1 \rightarrow$ Independent variable (IV), $b_1 \rightarrow$ Coefficient of x_1

Example - How a person salary is dependent on Experience.



$$\text{Salary} = b_0 + b_1 (\text{Experience})$$

$b_0 \rightarrow$ It is constant where the line crosses the vertical axis. = 30K

Inference, when Experience is 0, then the average starting salary is ₹30K. i.e., Fresher salary.

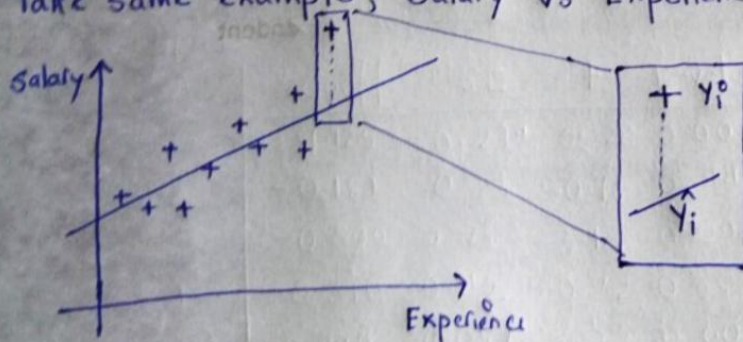
$b_1 \rightarrow$ the value of slope in unit change of x / Experience. Suppose we want to see how much salary change is there if experience is increased by one year. Inference, if the experience is increased from 3 to 4, then salary increased by ₹10K i.e., 50K to 60K.

So if we want to build the eqn, $\text{Salary} = 30,000 + b_1 (\text{Experience})$

$b_1 \neq 10K$ because that is change in y , but b_1 is the value of slope.

One of the method to find best fit line, is **Ordinary Least Square (OLS)**

Take same example, Salary vs Experience.



$y_i^o \rightarrow$ Actual Salary.

$\hat{y}_i \rightarrow$ Predicted Salary

$(y_i^o - \hat{y}_i)^2 \rightarrow$ Error.

Then for this regression line, we find the **total error**. $\rightarrow \sum (y - \hat{y})^2$

\rightarrow So for this regression line, we will calculate the total error. $\sum_1 (y - \hat{y})^2$.
 \rightarrow We will draw another line, and for that line we will calculate some $\sum_2 (y - \hat{y})^2$.
 \rightarrow We will choose the minimum $\sum_n (y - \hat{y})^2$, which will be the best line to find fit.

This method is known as Ordinary Least Square.

\rightarrow Multiple Linear Regression is on next page.

Simple linear regression,
Multiple linear regression,

$$y = b_0 + b_1 x_1 + b_2 x_2 + \dots + b_n x_n$$

Dependent variable (DV) y
Independent variables (IV) x_1, x_2, \dots, x_n
Constant b_0
Coefficients b_1, b_2, \dots, b_n

Assumption of linear regression

- i) Linearity
- ii) Homoscedasticity
- iii) Multivariate normality
- iv) Independence of errors
- v) Lack of multicollinearity

Dummy Variables

→ Suppose we have a categorical column which consist of set of alphabet values then we need to change the alphabets to numeric. That conversion is known as dummy variable. **Dummy Variable Trap**

Dependent Variable

Profit	Marketing	state/city
2000	1000	Pune
3000	1500	Mumbai
4000	1300	Mumbai
2500	1700	Pune

Dummy Variables

Profit	Marketing	Pune	Mumbai
2000	1000	1	0
3000	1500	0	1
4000	1300	0	1
2500	1700	1	0

$$(D_2 = 1 - D_1)$$

Dummy Variable Trap

— Pune & Mumbai columns are mirror image / switches. If one column is 1, other will be 0. So we can remove one column. This issue can be raised under multicollinearity. And in assumption, we said there should be no multicollinearity.

5 methods of building regression models

1. All in variable

2. Backward elimination

3. Forward selection

4. Bidirectional Elimination

5. Score comparison

Stepwise Regression

Some only bidirectional elimination is only refer as stepwise regression.

1. All in → Choose all the variable first and start removing wasteful variable.

2. Backward elimination → ① select a significance level to stay in model ($SL = 0.05$)
② Fit the model with all predictor variables.
③ Consider the predictor with highest P value. If $P > SL$, go ④
④ Remove the predictor.
⑤ Fit the model without this variables. (Final model)

3. Forward Selection → ① Select a significance level to enter the model ($SL = 0.05$)
② Fit all simple regression model $y \sim x_n$, select one with lowest P value.
③ Keep this variable and fit all possible models with one extra predictor added to one(s) already have.
④ Consider the predictor with lowest P-value. (Final model)

4. Bidirectional elimination → ① Select a significance level to enter and stay in a model.
② Perform the next step of forward selection ③ Perform all backward elimination steps
④ No new variable can enter (select) or delete (backward) Final model

5. Score Comparison — Adjusted R square, P value, AIC, BIC.